

Francisco Boudagh (13 h)

The three folders containing ema

1. Preproce

1.1 Look at a few emails from easy_ham, hard_ham and spam

```
import matplotlib.p
import seaborn as s
```

```
import numpy as np
import os
```

[illegible][illegible]

1.2 Note that

```
from sklearn.model_selection import train_test_split

if __name__ == '__main__':
    # Extract emails from 'easy han' and 'ohan'
    emails = extract_emails('easy han', 'ohan')
```

```
df_combined_eh_s = pd.concat([df_eh, df_s])

X = df_combined_eh_s['content']
y = df_combined_eh_s['class']

# splitting training data, train size: 65%, test size: 35%
hantrain, hanntest, spantrain, spanntest = train_test_split(X, y, train_size=0.65)
```

```
hamtrain = pd.DataFrame({'content': hamtrain, 'class': 'hamtrain'})
hamtest = pd.DataFrame({'content': hamtrain, 'class': 'hamtest'})
spamtrain = pd.DataFrame({'content': spamtrain, 'class': 'spamtrain'})
spamttest = pd.DataFrame({'content': spamttest, 'class': 'spamttest'})
```

2.1 Write a Python program that:

1. Uses the four datasets from Question 1 (`hamtrain`, `spamtrain`, `hamtest`, and `spamttest`)
2. Trains a Naive Bayes classifier (use the [scikit-learn library](#)) on `hamtrain` and `spamtrain`, that classifies the test sets and returns True Positive and False Negative rates on the `hamtest` and `spamttest` datasets. Use `CountVectorizer` ([Documentation here](#)) to transform the email texts into vectors. Please note that there are different types of Naive Bayes Classifier in scikit-learn ([Documentation here](#)). Test two of these classifiers that are well suited for this problem:
 - Multinomial Naive Bayes
 - Bernoulli Naive Bayes.

```
In [11]: from sklearn.feature_extraction.text import CountVectorizer
from sklearn.naive_bayes import MultinomialNB
from sklearn.metrics import confusion_matrix
from sklearn.metrics import classification_report
```

```
# cm labels
categories= ["Span", "Ham"]
```

- ```
def spam_detection_cm(hamtrain, hamtest, spamtrain, spamest, categories, groupnames, cm_title, classifier):
 vectorizer = CountVectorizer()
 hamtrain = vectorizer.fit_transform(hamtrain)
 hamtest = vectorizer.transform(hamtest)

 if classifier == 'RNB':
 model = MultinomialNB()
 elif classifier == 'nnb':
```

```
else:
 raise ValueError
```

```
fit train data into the chosen model
model.fit(hamtrain, spamtrain)

accuracy score in %
accuracy = model.score(hamtest, spamtest) * 100
print("Accuracy of this model is: {:.2f}%".format(accuracy))

spam_test_pred = model.predict(hamtest)
cm = confusion_matrix(spamtest, spam_test_pred)

group_counts = ["{0:0.0f}".format(value) for value in cm.flatten()]
group_percentages = ["{0:.2%}".format(value) for value in cm.flatten()/np.sum(cm)]
df_cm = pd.DataFrame(cm, range(2), range(2))

sns.set_style("whitegrid", {'axes.grid': False})
ax = sns.heatmap(df_cm, annot=True, fmt='', cmap='Blues', xticklabels=categories, yticklabels=categories)
ax.set_title(cm.title)
ax.set_xlabel("True")
ax.set_ylabel("Predicted")
plt.show()
```

```
span = 1 and ham = 0
df_ham_span.loc[df_ham_span['class']]
```

- ```

# create datasets with 65% train size and 35% test size
hamtrain, hamtest, spamtrain, spamest = train_test_split(df_ham_spam['content'], df_ham_spam['class'].values.tolist(), train_size=0.65)

spam_detection_cm(hamtrain, hamtest, spamtrain, spamest, categories, group_names, "Bernoulli Naive Bayes (easy ham)", "MNB")
spam_detection_cm(hamtrain, hamtest, spamtrain, spamest, categories, group_names, "Bernoulli Naive Bayes (easy ham)", "BNB")

```
- Accuracy of this model is: 97.01%
- Multinomial Naive Bayes (easy ham)

Confusion matrix for Bernoulli Naïve Bayes (easy ham):

	Predicted Spam	Predicted Ham
True Spam	128	31
True Ham	1	909

Accuracy of this model is: 92.14%

81

	Spam	Ham
True	6	904

Run (don't retrain) the two models from Question 2

- ```
df_ham_span = pd.concat([df_rh, df_sl])

df_span = 1 and ham = 0
df_ham_span.loc[df_ham_span['class'] == 'spam', 'class'] = 1
df_ham_span.loc[df_ham_span['class'] == 'ham', 'class'] = 0

Create datasets with 85% train size and 15% test size
hamtrain, hamtest, spamtrain, spamtest = train_test_split(df_ham_span['content'], df_ham_span['class'].values.tolist(), train_size=0.85)
```

Accuracy of this model is: 0.99

Confusion matrix for Multinomial Naive Bayes (hard ham) model:

|           | Predicted Spam | Predicted Ham |
|-----------|----------------|---------------|
| True Spam | 170            | 7             |
| True Ham  | 15             | 71            |

Accuracy of this model is: 88.21%

Page 10 of 10

Page 10 of 10

|           | Predicted Spam | Predicted Ham |
|-----------|----------------|---------------|
| True Spam | 85             | 120           |
| True Ham  | 31             | 55            |

```
df_h = extr
```

```
df_han_span = pd.concat([df_h, df_s])
```

## df\_ham\_span.2

```
hamtrain, hamtest, spamtrain, spamtest = train_test_split(df_ham_spam['content'], df_ham_spam['class'], values.tolist(), train_size=0.65)

spam_detection_cm(hamtrain, hamtest, spamtrain, spamtest, categories, group_names, "Multinomial Naïve Bayes (easy+hard ham)", "NBH")
spam_detection_cm(hamtrain, hamtest, spamtrain, spamtest, categories, group_names, "Bernoulli Naïve Bayes (easy+hard ham)", "BNB")
```

Accuracy of this model is: 95.82%

| Category | Accuracy |
|----------|----------|
| ham      | 95.82%   |
| spam     | 95.82%   |

Confusion matrix for Bernoulli Naive Bayes (easy+hard ham):

|           | Predicted Spam | Predicted Ham |
|-----------|----------------|---------------|
| True Spam | 75             | 7             |
| True Ham  | 4              | 177           |

Accuracy of this model is: 88.97%

un

5

|      | Spam | Ham |
|------|------|-----|
| Ham  | 0    | 180 |
| Spam | 1    | 0   |